

Support Vector Machines (SVM)

SVMs combine all the methods and ideas seen so far in order to find the best solution and also be robust to outliers.

Among all the linear models seen so far, this is the practical one.
(concept learning)

Let's consider a binary classification problem $y: \mathbf{x} \rightarrow \{+1, -1\}$ with dataset $\mathcal{D} = \{(x_m, t_m)\}_{m=1}^N$, $t_m \in \{+1, -1\}$ and a linear model

Assume \mathcal{D} is linearly separable:

$$y(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + w_0$$

linearly separable: it exists an hyperplane characterized by \mathbf{w} and w_0 that perfectly partitions the dataset.

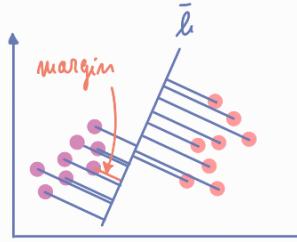
$$\exists (\mathbf{w}, w_0) \mid y(x_m) > 0 \text{ if } t_m = +1, y(x_m) < 0 \text{ if } t_m = -1$$

Given the geometric interpretation of a linear model, we want to maximize the distance between all the samples and the hyperplane.

MARGIN: smallest distance between a sample and the decision boundary.

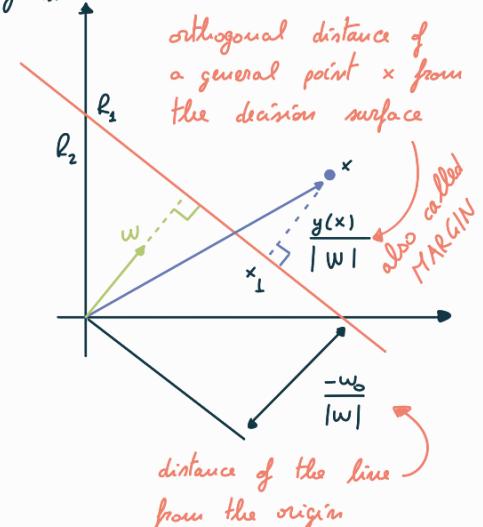
In SVMs the decision boundary is chosen to be the one for which the margin is maximized.

$$\text{margin} = \min_{m=1, \dots, N} \frac{|y(x_m)|}{\|\mathbf{w}\|} = \frac{1}{\|\mathbf{w}\|} \min_{m=1, \dots, N} |\mathbf{w}^T \mathbf{x}_m + w_0|$$



the margin is the minimum among the distances of all the samples in the dataset with respect to a particular hyperplane

$$\bar{t} = \bar{\mathbf{w}}^T \mathbf{x} + \bar{w}_0$$

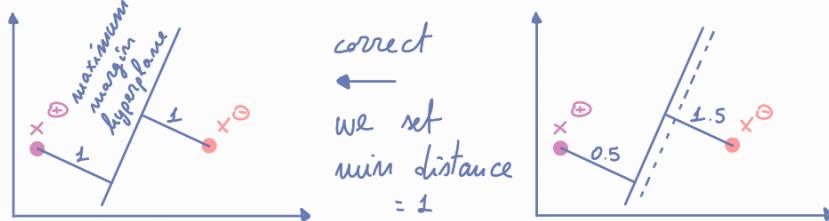


What we want: $\mathbf{w}^*, w_0^* = \underset{\mathbf{w}, w_0}{\operatorname{argmax}} \frac{1}{\|\mathbf{w}\|} \min_{m=1, \dots, N} |\mathbf{w}^T \mathbf{x}_m + w_0| \quad / \mathbf{w}, w_0 \text{ that maximize the margin}$

Direct optimization problem would be very complex to solve but we can convert it into an equivalent problem much easier to solve:

Rescale the dataset such that for the closest point x_k we have $t_k(\mathbf{w}^T \mathbf{x}_k + w_0) = 1$.
All the other points will have $t_m(\mathbf{w}^T \mathbf{x}_m + w_0) \geq 1$ ($\forall m = 1, \dots, N$).

When we find the optimal hyperplane, there will exist at least one point for each class such that its distance is exactly 1.



it cannot be other than that otherwise the chosen \bar{t} won't be the best hyperplane

from this formulation we can move to the lagrangian multipliers

We will use only these points and discard all the others.

$$\mathbf{w}^*, w_0^* = \underset{\mathbf{w}, w_0}{\operatorname{argmax}} \frac{1}{\|\mathbf{w}\|} = \underset{\mathbf{w}, w_0}{\operatorname{argmin}} \frac{1}{2} \|\mathbf{w}\|^2$$

This is an example of a quadratic programming problem in which we are trying to minimize a quadratic function subject to a set of linear inequality constraints.

Quadratic programming problems can be solved with the Lagrangian method. We obtain:

$$w^* = \sum_{m=1}^N \alpha_m^* t_m x_m$$

α_m^* = Lagrange multipliers

Optimization problem for determining w (dimension $|X|$) transformed in an optimization problem for determining α (dimension $|S|$) → efficient when $|X| \ll |S|$

dimension of the dataset (one Lagrange multiplier for each sample)

Why are we transforming an optimization problem in few dimensions (e.g. 2) into an optimization problem in thousand dimensions?

↳ when the value of α_m^* is 0 then the corresponding sample does not contribute to the solution.

& This optimization problem satisfies the KKT condition:

Karush-Kuhn-Tucker (KKT) condition:

for each $x_m \in S$, either $\alpha_m^* = 0$ or $t_m y(x_m) = 1$, thus $t_m y(x_m) > 1$ implies $\alpha_m^* = 0$

if $t_m y(x_m) = 1$
then $\alpha_m^* \neq 0$

all the points that are
not the closest one

The points that are not the closest ones (one for each class) does not contribute

Support vectors: x_k such that $t_k y(x_k) = 1$ and $\alpha_k^* > 0$

$$SV = \{x_k \in S \mid t_k y(x_k) = 1\}$$

The solution only depends on the samples that are in this set

$$y(x) = \sum_{x_j \in SV} \alpha_j^* t_j x_j^T x + w_0^* = 0$$

$$\begin{aligned} y(\tilde{x}) &= w^* \tilde{x} + w_0 \\ w^* &= \sum_{x_j \in SV} \alpha_j^* t_j x_j \end{aligned}$$

Other vectors $x_m \notin SV$ do not contribute ($\alpha_m^* = 0$)

To compute w_0 :

$$y(\tilde{x}) = w_0 + \sum_{x_j \in SV} \alpha_j^* t_j \tilde{x}_j^T x_j$$

① Support vector $x_m \in SV$ satisfies $t_k y(x_k) = 1$

$$t_k \left(\sum_{x_j \in SV} \alpha_j^* t_j x_j^T x_k + w_0^* \right) = 1$$

② Multiply by t_k and using $t_k^2 = 1$:

$$w_0^* = t_k - \sum_{x_j \in SV} \alpha_j^* t_j x_k^T x_j$$

$$y(x) = w^T x + w_0 \rightarrow w_0 = y(x) - w^T x$$

Instead of using one particular support vector x_m to determine w_0 , a more stable solution is obtained by averaging over all the support vectors

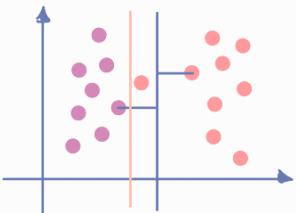
$$w_0^* = \frac{1}{|SV|} \sum_{x_k \in SV} \left(t_k - \sum_{x_j \in SV} \alpha_j^* t_j x_k^T x_j \right)$$

Given the maximum margin hyperplane determined by $\alpha_m^* w_0^*$, to classify a new instance:

$$\text{sign}(y(x')) = \text{sgn}\left(\sum_{x_k \in SV} \alpha_k^* t_k x'^T x_k + w_0^*\right)$$

we only need to determine the lagrange multipliers for the support vectors.

What if the data is almost linearly separable (e.g. few data points are on the other side of the true hyperplane)?



- hyperplane that the SVM we have seen so far would have found (worse)
- true hyperplane robust to noise

Even if ■ is the optimal solution for the actual dataset (accounting for noise) we still have that ■ will be better when we have to classify a new instance

How to deal with this? How to extend the problem in such a way that it can cope with this situation?

The constraint $t_m(w^T x_m + w_0) \geq 1$ means that all the points must be outside of the margin.

$$t_m y(x_m) \geq 1$$

We want to relax this constraint to allow some points to be inside the margin. We will introduce some variables in the formulation of the problem that will measure how much the constraint is violated: these variables will be called slack variables.

$$\varepsilon_m \geq 0 \quad m = 1, \dots, N$$

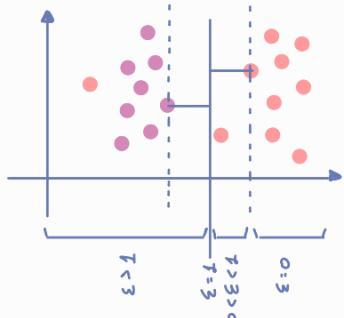
The new, less strict constraint will be called soft margin constraint.

$\varepsilon_m = 0$ if the point is inside the correct margin boundary

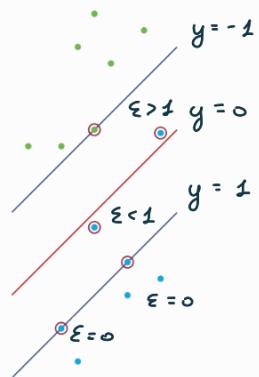
$0 < \varepsilon_m < 1$ if the point is inside the margin but still on the right side with respect to the optimal solution

$\varepsilon_m = 1$ if the point is exactly on the optimal plane

$\varepsilon_m > 1$ if the point is beyond the optimal plane



$\varepsilon > 1$ means that the sample will be misclassified by the solution



How to soften the constraint?

$$t_m y(x_m) \geq 1 \longrightarrow t_m y(x_m) \geq 1 - \varepsilon_m$$

We also want to minimize all the ε_m

$$w^*, w_0^* = \underset{w, w_0}{\text{argmin}} \left(\frac{1}{2} \|w\|^2 + C \sum_{m=1}^N \varepsilon_m \right)$$

first term to minimize is the maximum margin; second term to minimize is the sum of the slack variables

non-linearly separable data due to noise

C is a factor that weights how much you want to penalize the slack variables with respect to the maximization of the margin

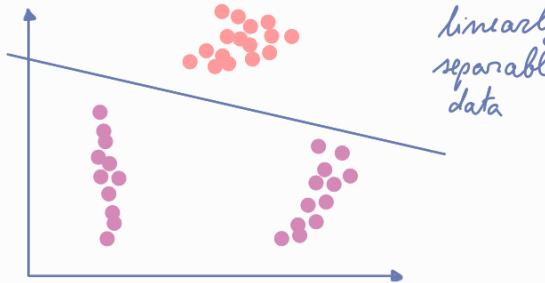
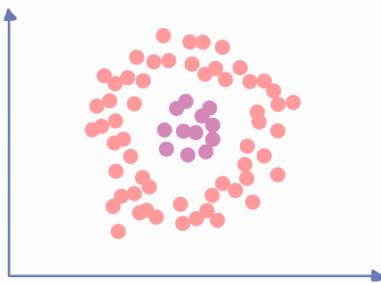
The solution of this optimization problem is similar to the previous one (this is still a quadratic programming problem that can be solved with the lagrangian optimization)

$$\omega^* = \sum_{m=1}^N \alpha_m^* t_m x_m \quad ; \quad \omega_0^* = \frac{1}{|SV|} \sum_{x_k \in SV} \left(t_k - \sum_{x_j \in SV} \alpha_j^* t_j x_k^T x_j \right)$$

This model is robust to outliers and to noise.

How to deal with non-linearly separable data due to the distribution itself (not to noise)?

intrinsically
non-linearly
separable data.
We can make a
change of coordinates
(cartesian \rightarrow polar)
and consider radius
as a feature



So far we considered models working directly on x . All the results hold if we consider a non-linear transformation of the inputs $\phi(x)$ called **basis function**. Decision boundaries will be linear in the feature space ϕ and non-linear in the original space x .
Classes that are linearly separable in the feature space ϕ may not be separable in the input space x .

$$D = \{(x_m, t_m)\} \quad \xrightarrow{\text{ } x_m \in \mathbb{R}^d} \quad D_\phi = \{(\phi_m = \phi(x_m), t_m)\}$$

$\phi_m = \phi(x_m) = \begin{bmatrix} \phi_0(x) \\ \vdots \\ \phi_n(x) \end{bmatrix} \quad n+1$

fixed to 1

We don't need to invert ϕ

SVM

- ↳ linearly separable data
- ↳ almost linearly separable data
- ↳ non-linearly separable data

Support Vector Machines (SVM)

$$x \in \mathbb{R}^d$$

We have a set of points we want to classify.
Each point is represented by a feature vector $x \in \mathbb{R}^d$

$$\phi: \mathbb{R}^d \rightarrow \mathbb{R}^n$$

$$\phi(x) \in \mathbb{R}^n$$

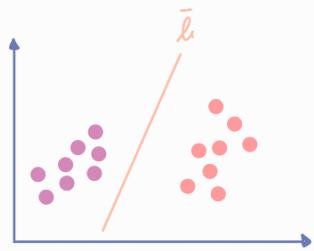
basis vector

This may be too simple for many applications,
so we want to map this to a more complex, non-linear
feature space ϕ .

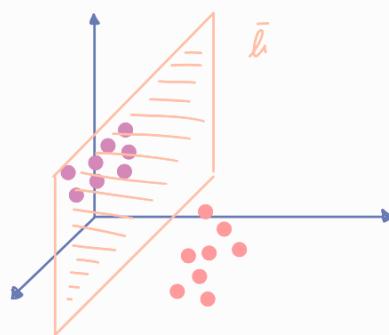
A decision boundary is a separator that divides these points into their respective classes.
This hyperplane is represented as:

$$H: w^T \phi(x) + w_0 = 0$$

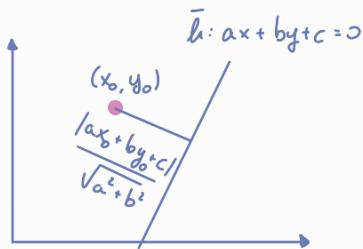
In a d -dimensional space our hyperplane is a $d-1$ -dimensional structure



in a 2-dimensional space the hyperplane is a 1-dimension line

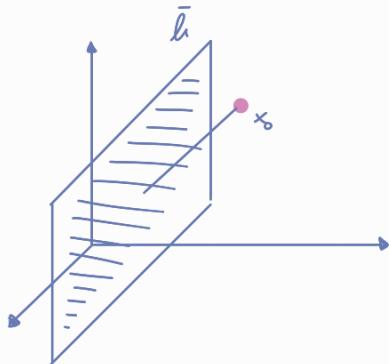


in a 3-dimensional space the hyperplane is a 2-dimension plane



the distance of a point from a line is this:

$$\frac{|ax_0 + by_0 + c|}{\sqrt{a^2 + b^2}}$$

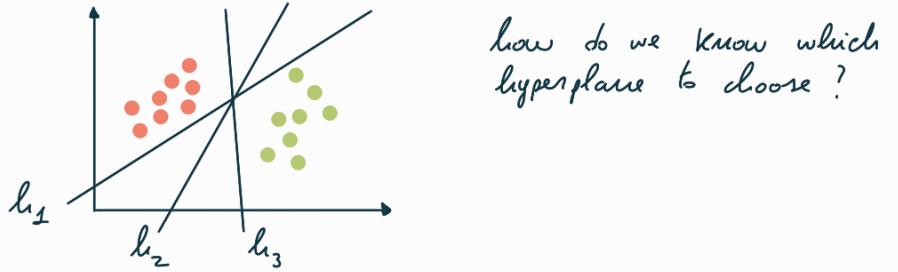


in the same way we can express the distance of a point from an hyperplane:

$$d_H(\phi(x_0)) = \frac{|w^T \phi(x_0) + w_0|}{\|w\|}$$

Consider the case where the data is perfectly separable. In this case it exists one hyperplane that can separate the training data groups with 100% accuracy.

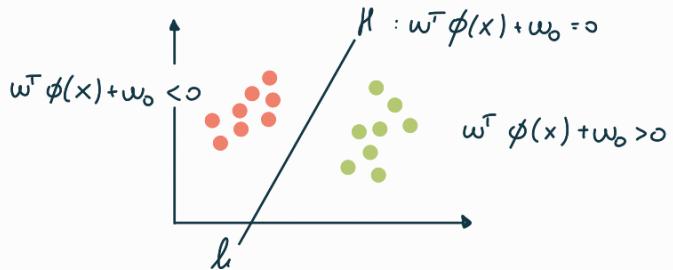
Consider a concept learning problem (classification of two classes):



Intuitively we can choose the hyperplane for which the distance from the closest points is maximum. This way we will make less mistakes during testing.

$$w^* = \operatorname{argmax}_w (\min_m d_H(\phi(x_m)))$$

The goal of SVM is defining the hyperplane that has the maximum margin from the closest points (in other words: maximize the minimum distance).



Stick with the concept learning problem for now:

Use this convention: $y_m \begin{cases} +1 & \text{for } \bullet \\ -1 & \text{for } \circ \end{cases}$

This way: $y_m (w^T \phi(x_m) + w_0) \begin{cases} \geq 0 & \text{correct} \\ < 0 & \text{incorrect} \end{cases}$

y_m	$w^T \phi(x_m) + w_0$	
≥ 0	≥ 0	≥ 0
< 0	< 0	> 0

ground truth prediction

② Perfect classification

$$w^* = \operatorname{argmax}_w (\min_m d_H(\phi(x_m))) = \operatorname{argmax}_w \left(\min_m \frac{|w^T \phi(x_m) + w_0|}{\|w\|} \right) \quad \left. \right\} \text{what we had so far}$$

$$= \operatorname{argmax}_w \left(\min_m \frac{|y_m (w^T \phi(x_m) + w_0)|}{\|w\|} \right) \quad \left. \begin{array}{l} \text{equal to the modulo} \\ \text{if all the samples} \\ \text{are correctly classified} \end{array} \right\} \text{which the new convention}$$

$$= \operatorname{argmax}_w \frac{1}{\|w\|} \left(\min_m |y_m (w^T \phi(x_m) + w_0)| \right) \quad \left. \begin{array}{l} \text{represents the} \\ \text{distance of the} \\ \text{closest point to } H \end{array} \right.$$

Let $\min_m |y_m (w^T \phi(x_m) + w_0)| = z$; with this normalization we impose that the distance of the closest point must be 1

We obtain:

$$w^* = \operatorname{argmax}_w \frac{1}{\|w\|}$$

this is possible by multiplying some constant factor c to the weight vector w and to the bias w_0

$$\begin{matrix} w &\leftarrow c w \\ w_0 &\leftarrow c w_0 \end{matrix}$$

From this formulation we can move to an equivalent one:

$$w^* = \arg\max \frac{1}{\|w\|} = \arg\min \frac{1}{2} \|w\|^2$$

This is the primal formulation of the SVM for perfectly separable data and assuming that all the points are classified correctly; in particular the latter assumption is often violated in reality since also the training set is subject to noise.

This is an example of a quadratic programming problem in which we are trying to minimize a quadratic function subject to a set of linear inequality constraints.

② Non-perfect classification

Instead of trying to classify every single point correctly and risk overfitting we allow to make some mistakes by using slack variables.

$$y_m (\mathbf{w}^T \phi(\mathbf{x}_m) + w_0) \geq 0 \quad \forall m$$

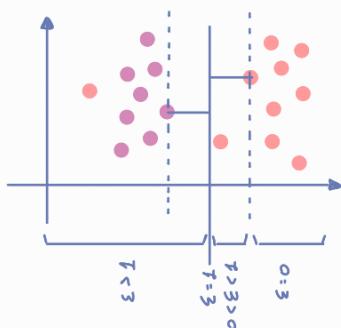
$$y_m (\mathbf{w}^T \phi(\mathbf{x}_m) + w_0) \leq 0 \quad \exists m$$

New primal form of SVM:

$$\min_{\mathbf{w}, w_0, \{\varepsilon_m\}} \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_m \varepsilon_m$$

$$y_m (\mathbf{w}^T \phi(\mathbf{x}_m) + w_0) \geq 1 - \varepsilon_m \quad \forall m \quad \varepsilon_m \geq 0 \quad \forall m$$

constraints



$\varepsilon > 1$ means that the sample will be missclassified by the solution

Low $C \rightarrow$ less complex boundary
High $C \rightarrow$ more complex boundary

These (case ① and ②) are examples of convex quadratic optimization problems: the objective function is quadratic in w and the constraints are linear in w and ε . To solve these constraint optimization problems we use the method of Lagrange multipliers.

Lagrange Multipliers

① Obtain the Primal form

$$\min_x f(x)$$

$$g_i(x) \leq 0$$

② Derive the Lagrangian function

$$L(x, \{\lambda_i\}) = f(x) + \sum_{i=1}^n \lambda_i g_i(x)$$

$\lambda_i \geq 0$ \downarrow Lagrange multipliers

③ Solve $\frac{\partial L}{\partial x} = 0$

④ We have a set of

$$x = h(\{\lambda_i\})$$

(isolate the original variables)

⑤ Substitute back in the Lagrangian

$$L(\{\lambda_i\}) = f(h(\{\lambda_i\})) + \sum_{i=1}^m \lambda_i g_i(h(\{\lambda_i\}))$$

⑥ Rewrite the problem

$$\max_{\{\lambda_i\} \geq 0} \min_x L(x, \{\lambda_i\})$$

We must have a function of only dual variables

Applied to our case we have:

$$\text{① } \min_{w, w_0, \{\varepsilon_m\}} \frac{1}{2} \|w\|^2 + c \sum_m \varepsilon_m \quad \left. \right\} \text{target function}$$

$$y_m (w^\top \phi(x_m) + w_0) \geq 1 - \varepsilon_m \quad \forall m \quad \left. \right\} \text{constraints}$$

$$\varepsilon_m \geq 0 \quad \left. \right\} \forall m$$

$$\textcircled{1} \rightarrow 1 - \varepsilon_m - y_m w^\top \phi(x_m) - y_m w_0 \leq 0$$

$$\textcircled{2} \rightarrow -\varepsilon_m \leq 0$$

② Determine the Lagrangian

$$L(w, w_0, \{\varepsilon_m\}, \{\lambda_m\}, \{\alpha_m\})$$

→ set of slack variables, one for each sample
 → set of lagrange multipliers, one for each sample for constraint ③
 → set of lagrange multipliers, one for each sample for constraint ②

$$\begin{aligned} &= f(h(\{\lambda_i\})) + \sum_{i=1}^m \lambda_i g_i(h(\{\lambda_i\})) \\ &= \frac{1}{2} \|w\|^2 + c \sum_m \varepsilon_m + \underbrace{\sum_m \alpha_m (1 - \varepsilon_m - y_m (w^\top \phi(x_m) + w_0))}_{\textcircled{1}} + \underbrace{\sum_m \lambda_m (-\varepsilon_m)}_{\textcircled{2}} \end{aligned}$$

③ Solve $\frac{\partial L}{\partial x} = 0$ (original variables are $w, w_0, \{\varepsilon_m\}$)

$$\frac{\partial L}{\partial w} = \frac{1}{2} \|w\|^2 + c \sum_m \varepsilon_m + \sum_m \alpha_m - \varepsilon_m \alpha_m - \alpha_m y_m w^\top \phi(x_m) - \alpha_m y_m w_0 - \sum_m \lambda_m \varepsilon_m = 0$$

$$\frac{\partial L}{\partial w_0} = \sum_m \alpha_m y_m = 0 \quad \text{recall these two rules for vector differentiation:}$$

$$\frac{\partial}{\partial x} (x^\top a) = \frac{\partial (a^\top x)}{\partial x} = a^\top \quad \text{with } x, a \in \mathbb{R}^{n \times 1}$$

$$\frac{\partial}{\partial \varepsilon_m} = c - \alpha_m - \lambda_m = 0$$

$$w - \sum_m \alpha_m y_m \phi(x_m) = 0$$

$$w = \sum_m \alpha_m y_m \phi(x_m)$$

② derivative of the transpose of a vector is the same as the derivative of the vector

⑤ Substitute these values back in the Lagrangian to eliminate the primal variables

$$\begin{aligned} L(\cdot) &= \frac{1}{2} \left(\sum_m \alpha_m y_m \phi(x_m) \right)^\top \left(\sum_m \alpha_m y_m \phi(x_m) \right) + c \sum_m \varepsilon_m + \sum_m (\alpha_m (1 - \varepsilon_m - y_m \left(\sum_m \alpha_m y_m \phi(x_m) \right)^\top \phi(x_m) \right. \\ &\quad \left. + w_0)) + \sum_m \lambda_m (-\varepsilon_m) \right) + \sum_m (\alpha_m (1 - \varepsilon_m - y_m \left(\sum_m \alpha_m y_m \phi(x_m) \right)^\top \phi(x_m) \right. \\ &\quad \left. + w_0)) + \sum_m \lambda_m (-\varepsilon_m) \end{aligned}$$

- old formula
- new formula

$$\begin{aligned}
 L(\cdot) &= \frac{1}{2} \sum_m \sum_n \alpha_m \alpha_n y_m y_n \phi^T(x_m) \phi(x_n) + c \sum_m \varepsilon_m + \sum_m \alpha_m - \sum_m \alpha_m \varepsilon_m \\
 &\quad - \sum_m \sum_n \alpha_m \alpha_n y_m y_n \phi^T(x_m) \phi(x_n) - w_0 \sum_m \alpha_m y_m - \sum_m \lambda_m \varepsilon_m \\
 &= -\frac{1}{2} \sum_m \sum_n \alpha_m \alpha_n y_m y_n \phi^T(x_m) \phi(x_n) + \dots \\
 \dots &= \sum_m \alpha_m - \frac{1}{2} \sum_m \sum_n \alpha_m \alpha_n y_m y_n \phi^T(x_m) \phi(x_n)
 \end{aligned}$$

this is the important part

This is the optimization problem that we want to solve:

$$\left\{
 \begin{array}{l}
 \max_{\alpha} \mathcal{O}(\alpha) = \sum_m \sum_n \alpha_m \alpha_n y_m y_n \phi^T(x_m) \phi(x_n) \\
 0 \leq \alpha_i \leq C \\
 \sum_m \alpha_m y_m = 0
 \end{array}
 \right.$$

This can be a very large QP optimization problem. We can use various techniques to tackle it but the best one is the Sequential Minimal Optimization Algorithm. I won't go into details.

The important part is that we have $\phi^T(x_m) \phi(x_n)$ and we will see that with the Kernelization we can compute it disregarding the true form of $\phi(\cdot)$.